

ITERATIVE METHODS FOR k -HESSIAN EQUATIONS

GERARD AWANOU

ABSTRACT. On a domain of the n -dimensional Euclidean space, and for an integer $k = 1, \dots, n$, the k -Hessian equations are fully nonlinear elliptic equations for $k > 1$ and consist of the Poisson equation for $k = 1$ and the Monge-Ampère equation for $k = n$. We prove a quadratic convergence rate for a finite difference approximation of smooth solutions of k -Hessian equations and the convergence of Newton's method. In addition we propose new iterative methods which are numerically shown to work for non smooth solutions. A connection of the latter with a popular Gauss-Seidel method for the Monge-Ampère equation is established and new Gauss-Seidel type iterative methods for 2-Hessian equations are introduced.

1. INTRODUCTION

Let Ω be a bounded, connected open subset of \mathbb{R}^n with boundary denoted $\partial\Omega$. Let $u \in C^2(\Omega)$ and for $x \in \Omega$, let $D^2u(x) = \left((\partial^2 u(x))(\partial x_i \partial x_j) \right)_{i,j=1,\dots,n}$ denote its Hessian. We denote the eigenvalues of $D^2u(x)$ by $\lambda_i(x), i = 1, \dots, n$. For $1 \leq k \leq n$, the k -Hessian operator is defined as

$$S_k(D^2u) = \sum_{i_1 < \dots < i_k} \lambda_{i_1} \cdots \lambda_{i_k}.$$

We note that $S_1(D^2u) = \Delta u$ is the Laplacian operator and $S_n(D^2u) = \det D^2u$ is the Monge-Ampère operator. We are interested in the numerical approximation of solutions of the Dirichlet problem for the k -Hessian equation

$$(1.1) \quad S_k(D^2u) = f \text{ in } \Omega, u = g \text{ on } \partial\Omega,$$

with f and g given and $f \geq 0$.

1.1. Quadratic convergence rate for a finite difference discretization. Let u^0 be a sufficiently close initial guess to the smooth solution u of (1.1). Consider the iterative method

$$(1.2) \quad \begin{aligned} \operatorname{div} \left(\{S_k^{ij}(D^2u^0)\} Du^{m+1} \right) &= \operatorname{div} \left(\{S_k^{ij}(D^2u^0)\} Du^m \right) + f - S_k(D^2u^m) \text{ in } \Omega \\ u^{m+1} &= g \text{ on } \partial\Omega, \end{aligned}$$

where $\{S_k^{ij}(D^2u^0)\}$ is a matrix which generalizes the cofactor matrix of D^2u^0 .

We prove the convergence of (1.2) at the continuous level in Hölder spaces and a quadratic convergence rate for the standard finite difference discretization. The quadratic convergence rate in the case of the Monge-Ampère equation was only known as "formally second order accurate" [3].

1.2. Newton's method. If one is only interested in smooth solutions, Newton's method is the most appropriate method. We analyze the convergence of Newton's method for solving (1.1) when it has a smooth solution.

1.3. Numerical work for subharmonicity preserving iterations. It is also of interest in some applications to be able to handle (1.1) when it has a non smooth solution. It has only been recently understood [1] that what is needed is a numerical method provably convergent for smooth solutions and numerically robust to handle non smooth solutions. The approach in [1] is to regularize the data and use approximation by smooth functions. The key to numerically handle non smooth solutions of (1.1) is to preserve convexity in the iterations. For discrete convexity we simply require discrete analogues of the usual notion of convexity for a natural discretization of D^2u using central difference approximations of the second order derivatives.

A smooth function u is said to be k -convex if $S_l(D^2u) \geq 0, 1 \leq l \leq k$. Convexity of a function can be shown to be equivalent to n -convexity, Lemma 2.4. Consider the iterative method

$$(1.3) \quad \Delta u^{m+1} = \left((\Delta u^m)^k + \frac{1}{c(k, n)} (f - S_k(D^2u^m)) \right)^{\frac{1}{k}} \text{ in } \Omega, u^{m+1} = g \text{ on } \partial\Omega,$$

with $c(k, n) = \binom{n}{k}/n^k$.

If D^2u has positive eigenvalues, we have the inequality

$$(1.4) \quad S_k(D^2u) \leq c(k, n)(\Delta u)^k,$$

which follows from the Maclaurin inequalities, [9, Proposition 1.1 (v i)].

For $k = 2$, (1.4) also holds with *no convexity assumption* on u , [16, Lemma 15.11]. Explicitly $c(2, 3) = 1/3$. Also, $c(n, n) = 1/n^n$ which gives

$$\det D^2u \leq \frac{1}{n^n} (\Delta u)^n,$$

a direct consequence of the arithmetic mean - geometric mean inequality.

If one starts with an initial guess u^0 such that $\Delta u^0 \geq 0$, (1.3) enforces $\Delta u^m \geq 0$ for all m . Indeed recall that $f \geq 0$ and assume that $\Delta u^m \geq 0$. Then by (1.4) $1/c(k, n)S_k(D^2u^m) \leq (\Delta u^m)^k$, and using (1.3) it follows that $(\Delta u^{m+1})^k \geq 0$. In other words, starting with an initial guess u_0 with $\Delta u^0 \geq 0$, (1.3) enforces subharmonicity.

Another class of iterative methods we introduce in this paper are Gauss-Seidel type iterative methods. The Gauss-Seidel methods are more efficient than (1.3) for large scale problems.

The simplicity of the methods discussed in this paper and the facility with which they can be implemented, make them attractive to researchers interested in Monge-Ampère equations. The other major motivation to study the subharmonicity preserving iterations is that they can be adapted to the finite element context and have been numerically shown in that context to be robust for non smooth solutions.

In two dimension (1.3) appears to perform well in the degenerate case $f \geq 0$ as discrete convexity is enforced in the iterations. The situation is different in three dimension. Numerical experiments indicate that in the degenerate case, with $k = 2, n = 3$, (1.3)

may not reproduce a solution which is convex. However, for $n = 3$ and $k = 3$, we can preserve convexity in the degenerate case by using the sequence of nonlinear 2-Hessian equations

$$(1.5) \quad S_2(D^2u^{m+1}) = 3 \left(\left(\frac{1}{3} S_2(D^2u^m) \right)^{\frac{3}{2}} + f - \det D^2u^m \right)^{\frac{2}{3}},$$

with $u_{m+1} = g$ on $\partial\Omega$. Each of these equations is solved iteratively by (1.3) with $k = 2, n = 3$. We note that $\left(\frac{1}{3} S_2(D^2u^m) \right)^{\frac{3}{2}} - \det D^2u^m \geq 0$ when $S_2(D^2u^m) > 0$, [16, Lemma 15.12]. Starting with an initial guess which satisfies $S_2(D^2u^0) > 0$ and setting $\det D^2u^m = 0$ in (1.5) whenever $S_2(D^2u^m) = 0$, we obtain a double sequence iterative method which at the limit enforces $\Delta u \geq 0$, $S_2 D^2u \geq 0$, and $\det D^2u = f \geq 0$. It was believed by many that preserving convexity in three dimension for the Monge-Ampère equation through a generalization of the iterative method introduced in [3] is not possible.

1.4. Relation with other work. The k -Hessian equations have mainly applications in conformal geometry and physics. The Monge-Ampère operator has received recently a lot of interest from numerical analysts. For $n = 3$ and $k = 2$, the numerical resolution of (1.1) has been considered in [18], where it was referred to as the σ_2 problem. The iterative method (1.3) generalizes an iterative method introduced in [3] for the two dimensional Monge-Ampère equation. The latter corresponds to the choice $k = n = 2$ and the constant $c(2, 2) = 1/4$ replaced by $1/2$. The 2-Hessian equation has also been considered recently in [17] from the point of view of monotone schemes.

We will see that if the finite difference discretization of (1.3) is solved by a Gauss-Seidel iterative method, one recovers a Gauss-Seidel iterative method which has been used by many authors to solve the two dimensional Monge-Ampère equation. We will refer to the latter method as the 2D Gauss-Seidel method for Monge-Ampère equation. It has been used in the numerical simulation of Ricci flow [11], as a smoother in multigrid methods for the balance vortex model in meteorology, [5, 4] and has been recently shown numerically to capture the viscosity solution of the 2D Monge-Ampère equation [3]. The connection between (1.3) and the 2D Gauss-Seidel method for Monge-Ampère equation is what enables us to introduce new Gauss-Seidel type iterative methods for k -Hessian equations. The ingredients of our proof of the quadratic convergence rate of the finite difference discretization are discrete Schauder estimates and a suitable generalization of the combined fixed point iterative method used in [8].

1.5. Organization of the paper. The paper is organized as follows: In the next section, we give some notations, recall the Schauder estimates and their discrete analogues. In section 3 we prove our main results on the quadratic convergence rate of a finite difference discretization of (1.1) and in section 4 we prove the convergence of Newton's method. In section 5 we introduce new Gauss-Seidel type iterative methods and their connections with the subharmonicity preserving iterations (1.3). Section 6 is devoted to numerical results. We conclude with some remarks.

Throughout the paper, for results at the continuous level, we will assume that the domain Ω is sufficiently smooth and uniformly convex. The reader interested only in the Monge-Ampère equation or for a first reading may assume that $k = n$.

Acknowledgments. The author is grateful to M. Neilan for many useful discussions. The author was supported in part by NSF grants DMS-0811052, DMS-1319640 and the Sloan Foundation.

2. NOTATION AND PRELIMINARIES

2.1. Hölder spaces and Schauder estimates. We denote by $C^r(\Omega)$ the set of all functions having all derivatives of order $\leq r$ continuous on Ω where r is a nonnegative integer or infinity and by $C^r(\overline{\Omega})$, the set of all functions in $C^r(\Omega)$ whose derivatives of order $\leq r$ have continuous extensions to $\overline{\Omega}$. For a multi-index $\beta = (\beta_1, \dots, \beta_n) \in \mathbb{N}^n$, put $|\beta| = \beta_1 + \dots + \beta_n$. We use the notation $D^\beta u(x)$ for the partial derivative $(\partial/\partial x_1)^{\beta_1} \dots (\partial/\partial x_n)^{\beta_n} u(x)$.

The norm in $C^r(\Omega)$ is given by

$$\|u\|_{r,\Omega} = \sum_{j=0}^r |u|_{j,\Omega}, \quad |u|_{j,\Omega} = \sup_{|\beta|=j} \sup_{\Omega} |D^\beta u(x)|.$$

A function u is said to be uniformly Hölder continuous with exponent α , $0 < \alpha \leq 1$ in Ω if the quantity

$$\sup_{x \neq y} \frac{|u(x) - u(y)|}{|x - y|^\alpha},$$

is finite. The space $C^{r,\alpha}(\overline{\Omega})$ consists of functions whose r -th order derivatives are uniformly Hölder continuous with exponent α in Ω . It is a Banach space with norm

$$\|u\|_{r,\alpha,\Omega} = \|u\|_{r,\Omega} + \sup_{|\beta|=r} \sup_{x \neq y} \frac{|D^\beta u(x) - D^\beta u(y)|}{|x - y|^\alpha}.$$

The norms $\|\cdot\|_{r,\Omega}$ and $\|\cdot\|_{r,\alpha,\Omega}$ are naturally extended to vector fields and matrix fields by taking the supremum over all components. We make the standard convention of using C for a generic constant. For $A = (a_{ij})_{i,j=1,\dots,n}$ and $B = (b_{ij})_{i,j=1,\dots,n}$ we recall that $A : B = \sum_{i,j=1}^n a_{ij} b_{ij}$. We will often use the following property

$$(2.1) \quad \|fg\|_{0,\alpha;\Omega} \leq C \|f\|_{0,\alpha;\Omega} \|g\|_{0,\alpha;\Omega}, \text{ for } f, g \in C^{0,\alpha}(\overline{\Omega}),$$

from which it follows that if A, B are matrix fields

$$(2.2) \quad \|A : B\|_{0,\alpha;\Omega} \leq C \sum_{i,j=1}^n \|a_{ij}\|_{0,\alpha;\Omega} \|b_{ij}\|_{0,\alpha;\Omega}.$$

We first state a global regularity result for the solution of strictly elliptic equations, which follows from [10, Theorems 6.14, 6.6 and Corollary 3.8].

Theorem 2.1. *Assume $0 < \alpha < 1$. Let Ω be a $C^{2,\alpha}$ domain in \mathbb{R}^n and $f \in C^\alpha(\overline{\Omega})$, $\phi \in C^{2,\alpha}(\overline{\Omega})$. We consider the strictly elliptic operator*

$$(2.3) \quad Lu = \sum_{i,j=1}^n a^{ij}(x) \frac{\partial^2}{\partial x_i \partial x_j} u(x),$$

with coefficients satisfying for positive constants λ, Λ ,

$$\sum_{i,j=1}^n a^{ij}(x) \zeta_i \zeta_j \geq \lambda \sum_{l=1}^n \zeta_l^2, \zeta_l \in \mathbb{R}, \text{ and } |a_{i,j}|_{0,\alpha;\Omega} \leq \Lambda.$$

Then the solution u of the equation

$$Lu = f \text{ in } \Omega, u = \phi \text{ on } \partial\Omega,$$

satisfies

$$\|u\|_{2,\alpha;\Omega} \leq C(\|\phi\|_{2,\alpha;\Omega} + \|f\|_{0,\alpha;\Omega}).$$

We will make the slight abuse of language of also denoting by $S_k(x), x = (x_1, \dots, x_n)$ the k th elementary symmetric polynomial of the variable x , i.e.

$$S_k(\lambda) = \sum_{i_1 < \dots < i_k} \lambda_{i_1} \cdots \lambda_{i_k}.$$

A function $u \in C^2(\Omega) \cap C^0(\overline{\Omega})$ with Hessian D^2u having eigenvalues $\lambda_i, i = 1, \dots, n$ is said to be k -admissible if $S_j(\lambda) > 0, j = 1, \dots, k$. Solutions of the k -Hessian equation will be required to be k -admissible, thus requiring $f > 0$. Moreover, let $\kappa = (\kappa_1, \dots, \kappa_{n-1})$ denote the principal curvatures of $\partial\Omega$. The domain Ω will be required to be $(k-1)$ -convex, i.e. there exists $c_0 > 0$ such that

$$S_j(\kappa) \geq c_0 > 0 \text{ on } \partial\Omega.$$

We then have, ([20, Theorems 3.3 and 3.4])

Theorem 2.2. *Assume that Ω is $(k-1)$ -convex, $\partial\Omega \in C^{3,1}$, $f \in C^{1,1}(\overline{\Omega})$, $\inf f > 0$, $g \in C^{3,1}(\overline{\Omega})$. Then there is a unique k -admissible solution $u \in C^{3,1}(\overline{\Omega})$ to the Dirichlet problem (1.1) with $0 < \alpha < 1$.*

We will need some identities for the k -Hessian operator $S_k(D^2u)$ which are derived explicitly for example in [9, p. 5–6]. See also [20]. For a symmetric matrix $A = (a_{ij})_{i,j=1,\dots,n}$ with eigenvalues $\lambda_i, i = 1, \dots, n$, let us also denote by $S_k(A)$ the k -th elementary symmetric polynomial of λ . This is equivalent to say that $S_k(A)$ is the sum of all $k \times k$ principal minors of A . Using the permutation definition of the determinant, we have

$$(2.4) \quad S_k(A) = \frac{1}{k!} \sum_{1 \leq i_1, \dots, i_k \leq n} \delta_{i_1, \dots, i_k}^{j_1, \dots, j_k} a_{i_1 j_1} \cdots a_{i_k j_k},$$

where $\delta_{i_1, \dots, i_k}^{j_1, \dots, j_k}$ is the generalized Kronecker delta which takes the value $+1$ if i_1, \dots, i_k differs from j_1, \dots, j_k by an even permutation and the value -1 otherwise. In other words, for a choice of i_1, \dots, i_k , $\delta_{i_1, \dots, i_k}^{j_1, \dots, j_k}$ is the signature of the permutation σ defined by $\sigma(i_l) = j_l, l = 1, \dots, k$. This implies that we only consider the case where the sets $\{i_1, \dots, i_k\}$ and $\{j_1, \dots, j_k\}$ are identical. Moreover we define $\delta_{i_1, \dots, i_k}^{j_1, \dots, j_k}$ to be 0 if $\{i_1, \dots, i_k\} \neq \{j_1, \dots, j_k\}$. Note also that $\{i_1, \dots, i_k\}$ is a subset of k elements of $\{1, \dots, n\}$.

We have

$$S_k^{ij}(A) := \frac{\partial}{\partial a_{ij}} S_k(A) = \frac{1}{(k-1)!} \sum_{1 \leq i, i_1, \dots, i_{k-1} \leq n} \delta_{i, i_1, \dots, i_{k-1}}^{j, j_1, \dots, j_{k-1}} a_{i_1 j_1} \cdots a_{i_{k-1} j_{k-1}},$$

and so $S_k(A) = \frac{1}{k} \sum_{i,j=1}^n S_k^{ij}(A) a_{ij}$. Let us denote by $\{S_k^{ij}(A)\}$ the symmetric matrix with entries $S_k^{ij}(A)$. We can write $S_k(A) = 1/k \{S_k^{ij}(A)\} : A$, that is $S_k(D^2v) = \frac{1}{k} \{S_k^{ij}(D^2v)\} : D^2v$. Using (2.4) and observing that the expression of $S_k(A)$ can be written in terms of a multilinear map, we obtain

$$(2.5) \quad S'_k(D^2v)D^2w = \{S_k^{ij}(D^2v)\} : D^2w.$$

Let us denote by $\{S_k^{ij}(A)\}'$ the Fréchet derivative of the mapping $A \rightarrow \{S_k^{ij}(A)\}$. Since $\{S_k^{ij}(A)\}'(B)$ is a sum of terms each of which is a product of $k-2$ terms from A and is linear in B , we have

$$(2.6) \quad \|\{S_k^{ij}(D^2v)\}'D^2w\|_{0,\Omega} \leq C|v|_{2,\Omega}^{k-2}|w|_{2,\Omega}.$$

Using (2.2) and (2.6) we also have

$$(2.7) \quad \|\{S_k^{ij}(D^2v)\}'D^2w\|_{0,\alpha;\Omega} \leq C|v|_{2,\alpha;\Omega}^{k-2}|w|_{2,\alpha;\Omega}.$$

Finally we note that

Lemma 2.3. *Let v be a strictly convex function with smallest eigenvalue uniformly bounded below by a constant $a > 0$. Then for $\eta = a/(2n)$, we have w strictly convex, whenever $\|w - v\|_{C^2(\Omega)} < \eta$.*

Proof. It follows from [12, Theorem 1 and Remark 2 p. 39] that for two symmetric $n \times n$ matrices A and B ,

$$(2.8) \quad |\lambda_l(A) - \lambda_l(B)| \leq n \max_{i,j} |A_{ij} - B_{ij}|, l = 1, \dots, n.$$

It follows that for $u, v \in C^2(\Omega)$,

$$(2.9) \quad |\lambda_1(D^2u(x)) - \lambda_1(D^2v(x))| \leq n\|w - v\|_{C^2(\Omega)}.$$

The result then follows. □

We conclude this section with the equivalence of n -convexity and convexity in the usual sense.

Lemma 2.4. *A C^2 function u is convex if and only if it is n -convex.*

Proof. If u is C^2 , $\lambda_i \geq 0$ on Ω for all i and thus $S_l(D^2u) \geq 0, l = 1, \dots, n$.

Conversely let us assume that A is a symmetric matrix with $S_l(A) \geq 0, l = 1, \dots, n$. We show that its eigenvalues λ_i are all positive. Let

$$p(\lambda) = \lambda^n + c_1\lambda^{n-1} + \dots + c_n,$$

denote the characteristic polynomial of A . It can be shown [13, Theorem 1.2.12] that

$$c_l = (-1)^l S_l(A), l = 1, \dots, n.$$

We show that if $\lambda_i < 0$ then $p(\lambda_i) \neq 0$. We have

$$\begin{aligned}
p(\lambda_i) &= \lambda_i^n + c_1 \lambda_i^{n-1} + \dots + c_n \\
&= \lambda_i^n + \sum_{l=1}^n (-1)^l S_l(A) \lambda_i^{n-l} \\
&= (-1)^n \left((-\lambda_i)^n + \sum_{l=1}^n (-1)^{l-n} S_l(A) \lambda_i^{n-l} \right) \\
&= (-1)^n \left((-\lambda_i)^n + \sum_{l=1}^n S_l(A) (-\lambda_i)^{n-l} \right).
\end{aligned}$$

Since $-\lambda_i > 0$ and $S_l(A) \geq 0$ for all l , we have $(-1)^n p(\lambda_i) \geq 0$. Moreover since $\sum_{l=1}^n S_l(A) (-\lambda_i)^{n-l} \geq 0$ and $-\lambda_i > 0$ we have $(-1)^n p(\lambda_i) \neq 0$. We conclude that $\lambda_i \geq 0$ for all i . This completes the proof. \square

2.2. Discrete Schauder estimates and related tools for the Poisson equation.

We will study the numerical approximation of (1.1)–(1.3) by standard finite difference discretizations. For simplicity, we consider a cuboidal domain Ω

$$\Omega = (0, 1)^n \subset \mathbb{R}^n, \text{ such that } \exists h, 0 < h < 1 \text{ with } 1/h \in \mathbb{Z}, i = 1, \dots, n.$$

Put

$$\begin{aligned}
\mathbb{Z}_h &= \{x = (x_1, \dots, x_n)^T \in \mathbb{R}^n : x_i/h \in \mathbb{Z}\} \\
\Omega_0^h &= \Omega \cap \mathbb{Z}_h, \Omega^h = \bar{\Omega} \cap \mathbb{Z}_h, \partial\Omega^h = \partial\Omega \cap \mathbb{Z}_h = \Omega^h \setminus \Omega_0^h.
\end{aligned}$$

Let $e^i, i = 1, \dots, n$ denote the i -th unit vector of \mathbb{R}^n . We define the following first order difference operators on the space $\mathcal{M}(\Omega^h)$ of grid functions $v^h(x), x \in \mathbb{Z}_h$,

$$\begin{aligned}
\partial_+^i v^h(x) &:= \frac{v^h(x + he^i) - v^h(x)}{h} \\
\partial_-^i v^h(x) &:= \frac{v^h(x) - v^h(x - he^i)}{h} \\
\partial_h^i v^h(x) &:= \frac{v^h(x + he^i) - v^h(x - he^i)}{2h}.
\end{aligned}$$

Higher order difference operators are obtained by combining the above difference operators. For a multi-index $\beta = (\beta_1, \dots, \beta_n) \in \mathbb{N}^n$, we define

$$\partial_+^\beta v^h := \partial_+^{\beta_1} \dots \partial_+^{\beta_n} v^h.$$

The operators ∂_-^β and ∂_h^β are defined similarly. Note that

$$(2.10) \quad \partial_+^i \partial_-^i v^h(x) = \frac{v^h(x + he^i) - 2v^h(x) + v^h(x - he^i)}{h^2}$$

$$\begin{aligned}
(2.11) \quad \partial_h^i \partial_h^j v^h(x) &= \frac{1}{4h^2} \left\{ v^h(x + he^i + he^j) + v^h(x - he^i - he^j) \right. \\
&\quad \left. - v^h(x + he^i - he^j) - v^h(x - he^i + he^j) \right\}, i \neq j.
\end{aligned}$$

The second order derivatives $\partial^2 v / \partial x_i \partial x_j$ are discretized using (2.10) and (2.11) for $i \neq j$. This gives a discretization of the Hessian $D^2 u$ which we denote by $\mathcal{H}_d(u^h)$.

Let v be a continuous function on Ω_h , we define $r_h(v)$ as the unique element of $\mathcal{M}(\Omega_h)$ characterized by

$$r_h(v)(x) = v(x), x \in \Omega_h,$$

and extend the operator r_h canonically to vector fields and matrix fields. For a function g defined on $\partial\Omega$, $r_h(g)$ defines the analogous restriction on $\partial\Omega_h$.

Thus the discrete version of (1.1) takes the form

$$(2.12) \quad S_k(\mathcal{H}_d u^h(x)) = r_h(f)(x), x \in \Omega^h, u^h(x) = r_h(g)(x) \text{ on } \partial\Omega.$$

The discrete Laplacian takes the form

$$(2.13) \quad \Delta_d(u^h) = \sum_{i=1}^n \partial_+^i \partial_-^i u^h,$$

while the discrete version of the linear elliptic operator (2.3) takes the form

$$L_d v^h(x) = \sum_{i,j=1}^n a^{ij}(x) \partial_h^i \partial_h^j v^h(x), x \in \Omega_h.$$

We now define discrete analogues of the Hölder norms and semi-norms following [14]. Let $[\xi, \eta]$ denote the set of points $\zeta \in \Omega^h$ such that $\xi_j \leq \zeta_j \leq \eta_j, j = 1, \dots, n$. Then for $v^h \in \mathcal{M}(\Omega^h), 0 < \alpha < 1$, we define

$$\begin{aligned} |v^h|_{j, \Omega_0^h} &= \max \{ |\partial_+^\beta v^h(\xi)|, |\beta| = j, [\xi, \xi + \beta] \subset \Omega^h \} \\ [v^h]_{j, \alpha, \Omega_0^h} &= \max \left\{ \frac{|\partial_+^\beta v^h(\xi) - \partial_+^\beta v^h(\eta)|}{(|\xi - \eta|)^\alpha}, |\beta| = j, \xi \neq \eta, [\xi, \xi + \beta] \cup [\eta, \eta + \beta] \subset \Omega^h \right\} \\ \|v^h\|_{p, \Omega_0^h} &= \max_{j \leq p} |v^h|_{j, \Omega_0^h} \\ \|v^h\|_{p, \alpha, \Omega_0^h} &= \|v^h\|_{p, \Omega_0^h} + [v^h]_{p, \alpha, \Omega_0^h}. \end{aligned}$$

The above norms are extended canonically to vector fields and matrix fields by taking the maximum over all components. For $j = 0$, we have discrete analogues of the maximum and $C^{0, \alpha}$ norms.

Note that Ω_0^h is an interior domain of Ω . We therefore have by [19, Theorem 2.1],

Theorem 2.5. *Assume $0 < \alpha < 1$ and $v^h = 0$ on $\partial\Omega^h$. Then there are constants C and h_0 such that for $v^h \in \mathcal{M}(\Omega^h), h \leq h_0$*

$$(2.14) \quad \|v^h\|_{2, \alpha, \Omega_0^h} \leq C(\|\Delta_d v^h\|_{0, \alpha, \Omega_0^h} + |v^h|_{0, \Omega_0^h}),$$

with the constant C independent of h .

Since

$$\begin{aligned} \partial_+^i \partial_-^i v^h(x) &= \partial_+^i \partial_+^i v^h(x - h e^i) \text{ and} \\ \partial_h^j \partial_h^i v^h(x) &= \frac{1}{4} \left(\partial_+^j \partial_+^i v^h(x) + \partial_+^j \partial_+^i v^h(x - h e^i) + \partial_+^j \partial_+^i v^h(x - h e^j) \right. \\ &\quad \left. + \partial_+^j \partial_+^i v^h(x - h e^i - h e^j) \right), \end{aligned}$$

we have $\max \{ \|\partial_+^i \partial_-^j v^h\|_{0,\alpha,\Omega_0^h}, \|\partial_h^j \partial_h^i v^h\|_{0,\alpha,\Omega_0^h}, i, j = 1, \dots, n \} \leq \|v^h\|_{2,\alpha,\Omega_0^h}$ and hence the above theorem also applies when the second order derivatives (2.10) and (2.11) are used in the definition of $\|\cdot\|_{2,\alpha,\Omega_0^h}$.

Next, we recall the discrete maximum principle which will allow us to control $|v^h|_{0,\Omega_0^h}$ in terms of $\|\Delta_d v^h\|_{0,\alpha,\Omega_0^h}$ and boundary terms. We have by [15, Theorem 2.1]

Theorem 2.6. *Assume that $\Delta_d v^h \geq r_h(f)$ in Ω_0^h . Then*

$$\max_{\Omega_0^h} v^h \leq \max_{\partial\Omega^h} \max\{v^h, 0\} + C(h) \left(\sum_{x \in \Omega_0^h} h^n |r_h(f)(x)|^n \right)^{\frac{1}{n}},$$

where $C(h) = c_1 h + c_2$, with c_1, c_2 positive constants independent of h and which depends only on n .

Since $\sum_{x \in \Omega_0^h} h^n \leq C$, using $-v^h$ in the above result and applying Theorem 2.5, we have the following discrete analogue of Theorem 2.1,

Theorem 2.7. *Assume that $\Delta_d v^h = f$ in Ω_0^h and $v^h = 0$ on $\partial\Omega^h$. Then there exists $h_1 > 0$ such that for $h \leq h_1$,*

$$\|v^h\|_{2,\alpha,\Omega_0^h} \leq (c_1 h + c_2) \|\Delta_d v^h\|_{0,\alpha,\Omega_0^h},$$

for positive constants c_1 and c_2 .

By Taylor series expansions, it is not difficult to verify that for $v \in C^2(\Omega)$

$$|r_h(v)|_{j,\Omega_0^h} \leq |v|_{2,\Omega}, j \leq 2.$$

Moreover, for $v \in C^4(\Omega)$,

$$(2.15) \quad \|r_h(D^2 v) - \mathcal{H}_d(r_h v)\|_{0,\alpha,\Omega_0^h} \leq C h^2 |v|_{4,\Omega},$$

and

$$[r_h(D^2 v) - \mathcal{H}_d(r_h v)]_{0,\alpha,\Omega_0^h} \leq C h^2 [v]_{4,\alpha,\Omega}.$$

To see that the last inequality holds, it is enough to consider a function of one variable $v \in C^4(-1, 1)$ and estimate $[v''(x) - (v(x+h) - 2v(x) + v(x-h))/h^2]_{0,\alpha}$. Now,

$$v''(x) - \frac{v(x+h) - 2v(x) + v(x-h)}{h^2} = \frac{h^2}{24} (v^{(4)}(x+t_1 h) + v^{(4)}(x-t_1 h)), t_1 \in [0, 1].$$

And for $t_1, t_2 \in [0, 1]$,

$$\frac{v^{(4)}(x+t_1 h) - v^{(4)}(y+t_2 h)}{|y-x|^\alpha} = \frac{v^{(4)}(x+t_1 h) - v^{(4)}(y+t_2 h)}{|y-x+h(t_2-t_1)|^\alpha} \frac{|y-x+h(t_2-t_1)|^\alpha}{|y-x|^\alpha}.$$

Recall that x and y grid points. Thus $y = x + sh$ for an integer s satisfying $1 \leq |s| \leq C$. This implies that $|y-x+h(t_2-t_1)|^\alpha/|y-x|^\alpha$ is bounded uniformly in x, y . The result then follows.

We have for $v \in C^{4,\alpha}(\Omega)$,

$$(2.16) \quad \|r_h(D^2 v) - \mathcal{H}_d(r_h v)\|_{0,\alpha,\Omega_0^h} \leq C h^2 \|v\|_{4,\alpha,\Omega}.$$

Lemma 2.8. *We have for $u \in C^{4,\alpha}(\Omega)$*

$$\|r_h(S_k D^2 u) - S_k \mathcal{H}_d(r_h u)\|_{0,\alpha,\Omega_h^h} \leq Ch^2 |u|_{2,\Omega}^{k-1} \|u\|_{4,\alpha,\Omega}.$$

Proof. By the mean value theorem, using (2.5), we have for some t in $[0, 1]$, and $x \in \Omega_0^h$,

$$\begin{aligned} S_k(D^2 u)(x) - S_k \mathcal{H}_d(r_h u)(x) &= S'_k(t D^2(u)(x) + (1-t) \mathcal{H}_d(r_h u)(x)) : (D^2 u(x) \\ &\quad - \mathcal{H}_d(r_h u)(x)) \\ &= \sum_{i,j=1}^n S_k^{ij}(t D^2(u)(x) + (1-t) \mathcal{H}_d(r_h u)(x)) (D^2 u(x) \\ &\quad - \mathcal{H}_d(r_h u)(x))_{ij}. \end{aligned}$$

Using (2.2), it follows that

$$\begin{aligned} \|r_h(S_k D^2 u) - S_k \mathcal{H}_d(r_h u)\|_{0,\alpha,\Omega_h^h} &\leq C(|u|_{2,\Omega} + |r_h u|_{2,\Omega_h^h})^{k-1} \|r_h(D^2 u) - \mathcal{H}_d(r_h u)\|_{0,\alpha,\Omega_h^h} \\ &\leq Ch^2 |u|_{2,\Omega}^{k-1} \|u\|_{4,\alpha,\Omega}. \end{aligned}$$

□

3. APPROXIMATIONS BY LINEAR ELLIPTIC PROBLEMS

In this section, we prove the convergence of the iterative method (1.2) and its discrete version. As indicated in the introduction, we also obtain the existence and uniqueness of the solution of the discrete version of (1.1), i.e. (2.12), as well as error estimates.

3.1. Convergence at the operator level. We assume that the assumptions of Theorem 2.2 hold. Then there is a unique k -admissible solution $u \in C^2(\Omega) \cap C^0(\bar{\Omega})$ of (1.1). Let $u_0 \in C^2(\Omega) \cap C^0(\bar{\Omega})$ such that $\|u - u_0\|_{2,\alpha,\Omega} < \delta$. For $k = n$, using an eigenvalue argument, it is not difficult to prove that the cofactor matrix is uniformly positive definite under the assumption $f \geq f_0 > 0$ for a constant f_0 . We assume that the matrix $\{S_k^{ij}(D^2 u)\}$ is uniformly positive definite. (The assumption is in fact true if one assumes that $S_l(D^2 u)$, $1 < l \leq k$ are uniformly bounded below. See [6] and [9, Theorem 1.3]. The complete justification is beyond the scope of this paper.)

By the continuity of the smallest eigenvalue of a matrix as a function of its entries, $\{S_k^{ij}(D^2 u^0)\}$ is also uniformly positive definite for $|u - u^0|_{2,\Omega}$ sufficiently small.

Next, $\{S_k^{ij}(D^2 u^0)\}$ is a symmetric matrix and divergence free by [9, Formula 1.10]. It is not difficult to check that for any sufficiently smooth matrix field A and vector field v , $\operatorname{div} A^T v = (\operatorname{div} A) \cdot v + A : Dv$ where the divergence of a matrix field is defined as the divergence operator applied row-wise. Thus we obtain

$$(3.1) \quad \operatorname{div} \left(\{S_k^{ij}(D^2 u^0)\} Dv \right) = \{S_k^{ij}(D^2 u^0)\} : D^2 v.$$

We have

Theorem 3.1. *Under the assumptions of Theorem 2.2, the sequence defined by (1.2) converges to u for u^0 sufficiently close to u .*

Proof. We define the operator $R : C^{2,\alpha}(\overline{\Omega}) \rightarrow C^{2,\alpha}(\overline{\Omega})$ by

$$\begin{aligned} -\operatorname{div} \left(\{S_k^{ij}(D^2 u^0)\} D(v - Rv) \right) &= -S_k(D^2 v) + f \text{ in } \Omega \\ R(v) &= g \text{ on } \partial\Omega. \end{aligned}$$

By Theorem 2.1, the operator R is well defined. We show that for $\rho > 0$ sufficiently small, R is a strict contraction in the ball $B_\rho(u) = \{v \in C^{2,\alpha}(\overline{\Omega}), \|u - v\|_{2,\alpha;\Omega} < \rho\}$.

For $v, w \in B_\rho(u)$ we have using (3.1)

$$\begin{aligned} \operatorname{div} \left(\{S_k^{ij}(D^2 u^0)\} D(Rv - Rw) \right) &= \operatorname{div} \left(\{S_k^{ij}(D^2 u^0)\} D(v - w) \right) + S_k(D^2 w) - S_k(D^2 v) \\ &= -\{S_k^{ij}(D^2 u^0)\} : (D^2 w - D^2 v) + S_k(D^2 w) - S_k(D^2 v). \end{aligned}$$

Next, by the mean value theorem and using (2.5), we have for some t in $[0, 1]$,

$$\begin{aligned} S_k(D^2 w) - S_k(D^2 v) &= \{S_k^{ij}(tD^2 w + (1-t)D^2 v)\} : D^2(w - v) \\ &= \{S_k^{ij}(t(D^2 w - D^2 u^0) + (1-t)(D^2 v - D^2 u^0) + D^2 u^0)\} : D^2(w - v). \end{aligned}$$

We use (2.7) to estimate the $C^{0,\alpha}$ norm of

$$A = \{S_k^{ij}(t(D^2 w - D^2 u^0) + (1-t)(D^2 v - D^2 u^0) + D^2 u^0)\} - \{S_k^{ij}(D^2 u^0)\}.$$

Put

$$\alpha_{st} = st(D^2 w - D^2 u^0) + s(1-t)(D^2 v - D^2 u^0) + D^2 u^0.$$

We have

$$(3.2) \quad |\alpha_{st}|_{0,\alpha;\Omega} \leq \|u_0 - v\|_{2,\alpha;\Omega} + \|u_0 - w\|_{2,\alpha;\Omega} + \|u_0\|_{2,\alpha;\Omega}.$$

By the mean value theorem, for some $s \in [0, 1]$ we have

$$A = \{S_k^{ij}(\alpha_{st})\}'(t(D^2 w - D^2 u^0) + (1-t)(D^2 v - D^2 u^0)),$$

and thus by (2.7)

$$(3.3) \quad \|A\|_{0,\alpha;\Omega} \leq C|\alpha_{st}|_{0,\alpha;\Omega}^{k-2}(\|u_0 - v\|_{2,\alpha;\Omega} + \|u_0 - w\|_{2,\alpha;\Omega}).$$

By Schauder estimates (Theorem 2.1), (2.2), (3.2) and (3.3) we obtain

$$\begin{aligned} (3.4) \quad \|R(v) - R(w)\|_{2,\alpha;\Omega} &\leq C\|A\|_{0,\alpha;\Omega}\|D^2(v - w)\|_{0,\alpha;\Omega} \\ &\leq C(\|u_0 - v\|_{2,\alpha;\Omega} + \|u_0 - w\|_{2,\alpha;\Omega} + \|u_0\|_{2,\alpha;\Omega})^{k-2} \\ &\quad (\|u_0 - v\|_{2,\alpha;\Omega} + \|u_0 - w\|_{2,\alpha;\Omega})\|v - w\|_{2,\alpha;\Omega} \\ &\leq C(\rho + \delta + \|u_0\|_{2,\alpha;\Omega})^{k-2}(\rho + \delta)\|v - w\|_{2,\alpha;\Omega}. \end{aligned}$$

Thus, for ρ and δ sufficiently small, R is a strict contraction.

It remains to show that R maps $B_\rho(u)$ into itself. We note by the definition of R and unicity of the solution of (1.1), a fixed point of R solves (1.1). Let $v \in B_\rho(u)$,

$$\|u - Rv\|_{2,\alpha;\Omega} = \|Ru - Rv\|_{2,\alpha;\Omega} \leq \|u - v\|_{2,\alpha;\Omega} \leq \rho,$$

which shows that R maps $B_\rho(u)$ into itself. The existence of a fixed point follows from the Banach fixed point theorem. Moreover, the sequence defined by $u^{m+1} = R(u^m)$, i.e. the sequence defined by (1.2), converges for ρ and δ sufficiently small to u . \square

3.2. Finite difference discretization. Next, we consider the following discrete version of (1.2)

(3.5)

$$\begin{aligned} \operatorname{div}_h \left(\{S_k^{ij}(\mathcal{H}_d u^{0,h})\} D_h u^{m+1,h} \right) &= \operatorname{div}_h \left(\{S_k^{ij}(\mathcal{H}_d u^{0,h})\} D_h u^{m,h} \right) \\ &\quad + f - S_k(\mathcal{H}_d u^{m,h}) \text{ in } \Omega \\ u^{m+1,h} &= g \text{ on } \partial\Omega, \end{aligned}$$

where for a mesh function v^h , we define $D_h v^h$ as the vector with components $\partial_-^i v^h$ and for a vector mesh function with components $v_i^h, i = 1, \dots, n$, $\operatorname{div}_h v^h = \sum_{i=1}^n \partial_+^i v_i^h$. With no requirement on the size of u , but otherwise under the assumptions of Theorem 3.5, we show that (2.12) has a unique solution to which the above sequence converges. Moreover, the convergence rate is $O(h^2)$. Define

$$(3.6) \quad B_\rho(r_h u) = \{v^h \in \mathcal{M}(\Omega^h), \|v^h - r_h u\|_{2,\alpha,\Omega_0^h} \leq \rho\}.$$

Lemma 3.2. *Let $S^h : \mathcal{M}(\Omega^h) \rightarrow \mathcal{M}(\Omega^h)$ be a strict contraction with contraction factor less than $1/2$, i.e. for $v^h, w^h \in \mathcal{M}(\Omega^h)$*

$$\|S^h(v^h) - S^h(w^h)\|_{2,\alpha,\Omega_0^h} \leq \frac{1}{2} \|v^h - w^h\|_{2,\alpha,\Omega_0^h}.$$

Let us also assume that S^h does not move the center $r_h(u)$ of the ball $B_\rho(r_h u)$ too far, i.e.

$$\|S^h(r_h u) - r_h u\|_{2,\alpha,\Omega_0^h} \leq C_0 h^2.$$

Then S^h maps $B_\rho(r_h u)$ into itself for $\rho = 2C_0 h^2$. Moreover S^h has a unique fixed point u_h in $B_\rho(r_h u)$ with the error estimate

$$\|r_h u - u^h\|_{2,\alpha,\Omega_0^h} \leq 2C_0 h^2.$$

Proof. For $v^h \in B_\rho(r_h u)$,

$$\begin{aligned} \|S^h(v^h) - r_h u\|_{2,\alpha,\Omega_0^h} &\leq \|S^h(v^h) - S^h(r_h u)\|_{2,\alpha,\Omega_0^h} + \|S^h(r_h u) - r_h u\|_{2,\alpha,\Omega_0^h} \\ &\leq \frac{1}{2} \|v^h - r_h u\|_{2,\alpha,\Omega_0^h} + C_0 h^2 \\ &\leq \frac{\rho}{2} + C_0 h^2 \leq \frac{\rho}{2} + \frac{\rho}{2} = \rho. \end{aligned}$$

This proves that S^h maps $B_\rho(r_h u)$ into itself. The existence of a fixed point follows from the Banach fixed point theorem. The convergence rate follows from the observation that

$$\begin{aligned} \|r_h u - u^h\|_{2,\alpha,\Omega_0^h} &\leq \|r_h u - S^h(r_h u)\|_{2,\alpha,\Omega_0^h} + \|S^h(r_h u) - S^h(u^h)\|_{2,\alpha,\Omega_0^h} \\ &\leq C_0 h^2 + \frac{1}{2} \|u^h - r_h u\|_{2,\alpha,\Omega_0^h}. \end{aligned}$$

□

Remark 3.3. *For h sufficiently small, $\mathcal{H}_d r_h(u)$ is sufficiently close to $D^2 u$ and hence $\{S_k^{ij}(\mathcal{H}_d r_h u)\}$ is positive definite, a property which also holds for $\{S_k^{ij}(\mathcal{H}_d u^{0,h})\}$ for $u^{0,h}$*

sufficiently close to $r_h(u)$. The arguments are similar to the ones of Lemma 2.3. See also Lemma 3.4 below.

We note

Lemma 3.4. *Let u be a k -admissible solution of (1.1). Assume that $\inf f > 0$ and $u \in C^4(\Omega)$. There exists $h_2 = C/\|u\|_{4,\Omega}$ such that for $h \leq h_2$, $\Delta_d(r_h u) \geq c_0 > 0$ where $c_0 = 1/2((\inf f)/c(k, n))^{1/k}$. Moreover, if u is a strictly convex function, then for $h \leq h_2$ and $\rho = O(h^2)$, $\mathcal{H}_d(r_h u)$ is a positive matrix and v^h is a discrete convex function, when $v^h \in B_\rho(r_h u)$.*

Proof. Since the eigenvalues of a matrix are continuous functions of its entries (as roots of the characteristic polynomial), for a matrix $A = (a_{ij})$ with $S_k A > 0$, we have for $\epsilon > 0$, the existence of $\gamma > 0$ depending only on the space dimension n such that $|S_k B - S_k A| < \epsilon$ when $\sup_{ij} |b_{ij} - a_{ij}| < \gamma$. This implies $S_k B > S_k A - \epsilon$. Thus with $\epsilon = (S_k A)/2$, we have $S_k B > (S_k A)/2$.

For $h \leq h_2 = C/\|u\|_{4,\Omega}$ we have $Ch^2\|u\|_{4,\Omega} < \gamma$ and thus since $S_k(D^2 u) = f > \inf f > 0$, by (2.15) $S_k(\mathcal{H}_d(r_h u)) \geq 1/2 \inf f$. This implies that $\mathcal{H}_d(r_h u)$ is a positive matrix. By (1.4)

$$\Delta_d(r_h u) \geq \frac{1}{2}((\inf f)/c(k, n))^{1/k}.$$

Let $v^h \in B_\rho(r_h u)$. Then by definition of $B_\rho(r_h u)$ and (2.15)

$$\begin{aligned} \|\mathcal{H}_d(v^h) - \mathcal{H}_d(r_h u)\|_{0,\alpha,\Omega_0^h} &\leq \|\mathcal{H}_d(v^h) - r_h(D^2 u)\|_{0,\alpha,\Omega_0^h} + \|r_h(D^2 u) - \mathcal{H}_d(r_h u)\|_{0,\alpha,\Omega_0^h} \\ &\leq \rho + Ch^2\|u\|_{4,\Omega}, \end{aligned}$$

which can be made smaller than γ for h and ρ sufficiently small. Thus given that $\mathcal{H}_d(r_h u)$ is positive definite, the same holds for $\mathcal{H}_d(v^h)$. \square

Theorem 3.5. *Assume $u \in C^{4,\alpha}(\Omega)$ and $\inf f > 0$. Choose $u^{0,h}$ such that $\|u^{0,h} - r_h u\|_{2,\alpha,\Omega_0^h} \leq \delta, \delta > 0$. For $h \leq h_2 = C/\|u\|_{4,\Omega}$ sufficiently small, (2.12) has a unique solution u^h which satisfies $\Delta_d(u^h) \geq 0$ and u^h converges to the unique solution u of (1.1) as $h \rightarrow 0$ with quadratic convergence rate for δ sufficiently small.*

Proof. We define the operator $R^h : \mathcal{M}(\Omega^h) \rightarrow \mathcal{M}(\Omega^h)$ by

$$\begin{aligned} -\operatorname{div}_h \left(\{S_k^{ij}(\mathcal{H}_d u^{0,h})\} D(v^h - R^h v^h) \right) &= -S_k(\mathcal{H}_d v^h) + f \text{ in } \Omega \\ R^h(v^h) &= v^h \text{ on } \partial\Omega, \end{aligned}$$

and show that R^h has a unique fixed point in $B_\rho(r_h u)$ for $\rho = O(h^2)$. By Remark 3.3 the above problem is then well defined.

Next, note that with (2.16) applied to u one has $|r_h(u)|_{2,\alpha,\Omega_0^h} \leq C\|u\|_{2,\alpha,\Omega}$. As in the proof of Theorem 3.1, see (3.4), R^h is a strict contraction in $B_\rho(r_h u)$ for ρ and δ sufficiently small. Moreover, the contraction factor can be made smaller than $1/2$ by choosing h and δ sufficiently small.

Since $f = S_k(D^2 u)$, by the discrete Schauder estimates (2.7) and Lemma 2.8

$$\|R^h(r_h u) - r_h u\|_{2,\alpha,\Omega_0^h} \leq C(c_1 h + c_2) \|r_h(S_k D^2 u) - S_k \mathcal{H}_d(r_h u)\|_{0,\alpha,\Omega_0^h} \leq Ch^2,$$

since $h \leq 1$. By Lemma 3.2 we conclude that R^h has a fixed point u^h in $B_\rho(r_h u)$ with the claimed convergence rate.

The claimed property of u^h follows from the fact that $u^h \in B_\rho(r_h u)$ and Lemma 3.4. \square

Remark 3.6. *As with Lemma 2.3, the constant δ which controls how close u^0 is to u , scales linearly with the size of u . Thus, if necessary, by rescaling the equation, i.e. solve $S_k(D^2\beta u) = \beta^k f$, for $\beta > 0$, it is always possible to find a suitable initial guess. Indeed let $\epsilon > 0$ denote a user's measure of the closeness of an initial guess. Since $\|u - u^0\|_{2,\alpha} \leq \delta$, we have $\|\beta u - \beta u^0\|_{2,\alpha} \leq \beta\delta$. One can therefore choose β such that $\beta\delta < \epsilon$.*

4. NEWTON'S METHOD

As in the previous section, we assume that $\{S_k^{ij}(D^2u)\}$ is uniformly positive definite. By Remark 3.3, for h sufficiently small, there exists $m' > 0$ such that for $v^h \in B_\rho(r_h u)$, $\{S_k^{ij}(\mathcal{H}_d v^h)\}$ has smallest eigenvalue greater than m' . We consider for $u^{0,h} \in B_\rho(r_h u)$ the sequence of iterates

$$(4.1) \quad \begin{aligned} \{S_k^{ij}(\mathcal{H}_d u^{m,h})\} : (\mathcal{H}_d u^{m+1,h} - \mathcal{H}_d u^{m,h}) &= r_h(f) - S_k(\mathcal{H}_d u^{m,h}) \text{ in } \Omega \\ u^{m+1,h} &= g \text{ in } \partial\Omega. \end{aligned}$$

Theorem 4.1. *The sequence defined by (4.1) satisfies*

$$(4.2) \quad \|u^{m+1,h} - u^h\|_{2,\alpha;\Omega} \leq C \|u^{m,h} - u^h\|_{2,\alpha;\Omega}^2,$$

for ρ and h sufficiently small and where u^h denotes the solution of (2.12) in $B_\rho(r_h u)$.

Proof. Put

$$(4.3) \quad B = \{S_k^{ij}(\mathcal{H}_d u^{m,h})\} : (\mathcal{H}_d u^{m+1,h} - \mathcal{H}_d u^h).$$

We have by (2.12)

$$(4.4) \quad \begin{aligned} B &= \{S_k^{ij}(\mathcal{H}_d u^{m,h})\} : (\mathcal{H}_d u^{m,h} - \mathcal{H}_d u^h) + S_k(\mathcal{H}_d u^h) - S_k(\mathcal{H}_d u^{m,h}) \\ &= \left(\{S_k^{ij}(\mathcal{H}_d u^{m,h})\} - \{S_k^{ij}(\mathcal{H}_d u^h)\} \right) : (\mathcal{H}_d u^{m,h} - \mathcal{H}_d u^h) \\ &\quad + \{S_k^{ij}(\mathcal{H}_d u^h)\} : (\mathcal{H}_d u^{m,h} - \mathcal{H}_d u^h) + S_k(\mathcal{H}_d u^h) - S_k(\mathcal{H}_d u^{m,h}). \end{aligned}$$

Put

$$(4.5) \quad B_1 = \left(\{S_k^{ij}(\mathcal{H}_d u^{m,h})\} - \{S_k^{ij}(\mathcal{H}_d u^h)\} \right) : (\mathcal{H}_d u^{m,h} - \mathcal{H}_d u^h),$$

and

$$(4.6) \quad B_2 = \{S_k^{ij}(\mathcal{H}_d u^h)\} : (\mathcal{H}_d u^{m,h} - \mathcal{H}_d u^h) + S_k(\mathcal{H}_d u^h) - S_k(\mathcal{H}_d u^{m,h}).$$

By the mean value theorem, (2.5) and (2.7), we have

$$B_1 = (\{S_k^{ij}(t\mathcal{H}_d u^{m,h} + (1-t)\mathcal{H}_d u^h)\}'(\mathcal{H}_d u^{m,h} - \mathcal{H}_d u^h)) : (\mathcal{H}_d u^{m,h} - \mathcal{H}_d u^h),$$

for $t \in [0, 1]$ and thus

$$\begin{aligned}
 \|B_1\|_{0,\alpha;\Omega} &\leq C(\|u^h\|_{2,\alpha;\Omega} + \|u^{m,h}\|_{2,\alpha;\Omega})^{k-2} \|u^{m,h} - u^h\|_{2,\alpha;\Omega}^2 \\
 (4.7) \quad &\leq C(\|r_h u\|_{2,\alpha;\Omega} + \rho)^{k-2} \|u^{m,h} - u^h\|_{2,\alpha;\Omega}^2 \\
 &\leq C(\|u\|_{2,\alpha;\Omega} + \rho)^{k-2} \|u^{m,h} - u^h\|_{2,\alpha;\Omega}^2.
 \end{aligned}$$

We also have by the mean value theorem

$$\begin{aligned}
 B_2 &= \{S_k^{ij}(\mathcal{H}_d u^h)\} : (\mathcal{H}_d u^{m,h} - \mathcal{H}_d u^h) \\
 &\quad + \{S_k^{ij}(t\mathcal{H}_d u^h + (1-t)\mathcal{H}_d u^{m,h})\} : (\mathcal{H}_d u^h - \mathcal{H}_d u^{m,h}) \\
 (4.8) \quad &= \left(\{S_k^{ij}(\mathcal{H}_d u^h)\} - \{S_k^{ij}(t\mathcal{H}_d u^h + (1-t)\mathcal{H}_d u^{m,h})\} \right) : (\mathcal{H}_d u^{m,h} - \mathcal{H}_d u^h) \\
 &= \left(\{S_k^{ij}((1-s)\mathcal{H}_d u^h + st\mathcal{H}_d u^h + s(1-t)\mathcal{H}_d u^{m,h})\}' \right. \\
 &\quad \left. ((1-t)(\mathcal{H}_d u^h - \mathcal{H}_d u^{m,h})) \right) : (\mathcal{H}_d u^{m,h} - \mathcal{H}_d u^h),
 \end{aligned}$$

for $s, t \in [0, 1]$. As for B_1 we obtain

$$(4.9) \quad \|B_2\|_{0,\alpha;\Omega} \leq C(\|u\|_{2,\alpha;\Omega} + \rho)^{k-2} \|u^{m,h} - u^h\|_{2,\alpha;\Omega}^2.$$

Combining (4.3)–(4.8) and using Schauder estimates, we obtain (4.2). Choosing ρ such that $C\rho < 1$, we conclude that $u^{m+1,h} \in B_\rho(r_h u)$ when $u^{m,h} \in B_\rho(r_h u)$ and the quadrate convergence rate of Newton's method. \square

Remark 4.2. *The proof of convergence of Newton's method given here can also be reproduced at the continuous level.*

5. GAUSS-SEIDEL ITERATIVE METHODS

It is a natural idea to solve (2.12) by a nonlinear Gauss-Seidel method, that is solve (2.12) for $u^h(x)$ and solve the resulting nonlinear equations by a Gauss-Seidel method. Although this seems a daunting task for arbitrary k , we show that for $k = 2$, this takes a very elegant form. We then establish a connection between the resulting nonlinear Gauss-Seidel iterative method for 2-Hessian equations and the discrete version of (1.3), i.e.

$$\begin{aligned}
 (5.1) \quad \Delta_d u^{m+1,h} &= \left((\Delta_d u^{m,h})^k + \frac{1}{c(k,n)} (r_h(f) - S_k(\mathcal{H}_d u^{m,h})) \right)^{\frac{1}{k}} \\
 u^{m+1,h} &= r_h(g) \text{ on } \partial\Omega,
 \end{aligned}$$

when the Gauss-Seidel method is used to solve the Poisson equations. To take advantage of that observation, we introduce a partial Gauss-Seidel iterative method for k -Hessian equations where updates are done only on the linear part of the scheme.

5.1. Nonlinear Gauss-Seidel method for 2-Hessian equations. We start with the identity

$$(5.2) \quad \Delta_d u^h = \left((\Delta_d u^h)^2 + \frac{1}{c(2, n)}(r_h(f) - S_2(\mathcal{H}_d u^h)) \right)^{\frac{1}{2}},$$

and show that the right hand side is independent of $u^h(x)$. Note that by (2.11), $\partial_h^i \partial_h^j v^h(x)$, $i \neq j$ is independent of $u^h(x)$ and by (2.13),

$$\frac{\partial(\Delta_d u^h(x))}{\partial(u^h(x))} = \sum_{i=1}^n -\frac{2}{h^2} = -\frac{2n}{h^2}.$$

Since $\partial S_k(A)/\partial z = \sum_{i,j=1}^n (\partial S_k(A)/\partial a_{ij})(\partial a_{ij}/\partial z)$, we conclude that

$$\begin{aligned} \frac{\partial}{\partial(u^h(x))} S_2(\mathcal{H}_d u^h(x)) &= \sum_{\substack{i,j=1 \\ i \neq j}}^n S_2^{ij}(\mathcal{H}_d u^h(x)) \frac{\partial}{\partial(u^h(x))} \partial_h^i \partial_h^j u^h(x) \\ &\quad + \sum_{i=1}^n S_2^{ii}(\mathcal{H}_d u^h(x)) \frac{\partial}{\partial(u^h(x))} \partial_+^i \partial_-^i u^h(x) \\ &= -\frac{2}{h^2} \sum_{i=1}^n S_2^{ii}(\mathcal{H}_d u^h(x)) = -\frac{2}{h^2} \sum_{i=1}^n \sum_{\substack{1 \leq p \leq n \\ p \neq i}} \delta_{ip}^p \partial_+^p \partial_-^p u^h(x) \\ &= -\frac{2}{h^2} \sum_{i=1}^n \sum_{p \neq i} \partial_+^p \partial_-^p u^h(x) = -\frac{2}{h^2} (n-1) \Delta_d u^h(x) \\ &= -\frac{2}{h^2} (2n) c(2, n) \Delta_d u^h(x) = -\frac{4n}{h^2} c(2, n) \Delta_d u^h(x), \end{aligned}$$

and we recall that the definition of δ_{ip}^{ip} was given in section 2.1. This gives

$$\frac{\partial}{\partial(u^h(x))} \left((\Delta_d u^h(x))^2 + \frac{1}{c(2, n)}(r_h(f) - S_2(\mathcal{H}_d u^h(x))) \right) = 0.$$

We can therefore rewrite (5.2) as

$$(5.3) \quad u^h(x) = \frac{h^2}{2n} \left[\sum_{i=1}^n \frac{u^h(x + he^i) + u^h(x - he^i)}{h^2} - \left((\Delta_d u^h(x))^2 + \frac{1}{c(2, n)}(r_h(f) - S_2(\mathcal{H}_d u^h(x))) \right)^{\frac{1}{2}} \right],$$

where the solution with $\Delta_d u^h \geq 0$ has been selected. For $n = 2$, this is the identity which was solved in [11, 5, 4, 3] by a Gauss-Seidel iterative method, as indicated in the introduction. For $n \geq 3$, this provides new iterative methods for the 2-Hessian equations.

Henceforth, we shall assume that a row ordering of the elements of Ω^h is chosen. Note that if we apply the Gauss-Seidel method to the problem (5.1), we obtain a double

sequence $u^{m,p,h}$ defined by

$$u^{m+1,p+1,h}(x) = \frac{h^2}{2n} \left[\sum_{i=1}^n \frac{u^{m+1,p,h}(x + he^i) + u^{m+1,p+1,h}(x - he^i)}{h^2} - \left((\Delta_d u^{m,h}(x))^2 + \frac{1}{c(2,n)}(r_h(f) - S_2(\mathcal{H}_d u^{m,h}(x))) \right)^{\frac{1}{2}} \right],$$

This leads us to consider the double sequence $u_p^{m,h}$ defined by

$$u_{p+1}^{m+1,h}(x) = \frac{h^2}{2n} \left[\sum_{i=1}^n \frac{u_p^{m+1,h}(x + he^i) + u_{p+1}^{m+1,h}(x - he^i)}{h^2} - \left((\Delta_d u_{p*}^{m,h}(x))^2 + \frac{1}{c(2,n)}(r_h(f) - S_2(\mathcal{H}_d u_{p*}^{m,h}(x))) \right)^{\frac{1}{2}} \right],$$

where $\Delta_d u_{p*}^{m,h}(x)$ and $S_2(\mathcal{H}_d u_{p*}^{m,h}(x))$ are the actions of the discrete Laplace and 2-Hessian operators on $u_p^{m,h}$ updated with the most recently computed values.

Formally, as $m \rightarrow \infty$, this gives the nonlinear Gauss-Seidel method

$$(5.4) \quad u_{p+1}^h(x) = \frac{h^2}{2n} \left[\sum_{i=1}^n \frac{u_p^h(x + he^i) + u_{p+1}^h(x - he^i)}{h^2} - \left((\Delta_d u_{p*}^h(x))^2 + \frac{1}{c(2,n)}(r_h(f) - S_2(\mathcal{H}_d u_{p*}^h(x))) \right)^{\frac{1}{2}} \right],$$

where as above $\Delta_d u_{p*}^h(x)$ and $S_2(\mathcal{H}_d u_{p*}^h(x))$ are the actions of the discrete Laplace and 2-Hessian operators on u_p^h updated with the most recently computed values of u_{p+1}^h .

5.2. Partial Gauss-Seidel method for 2-Hessian equations. One can also consider the following partial Gauss-Seidel iterative method where updates of u_p^h are done only on the linear part of the right hand side of (5.4):

$$(5.5) \quad u_{p+1}^h(x) = \frac{h^2}{2n} \left[\sum_{i=1}^n \frac{u_p^h(x + he^i) + u_{p+1}^h(x - he^i)}{h^2} - \left((\Delta_d u_p^h(x))^2 + \frac{1}{c(2,n)}(r_h(f) - S_2(\mathcal{H}_d u_p^h(x))) \right)^{\frac{1}{2}} \right].$$

6. NUMERICAL RESULTS

We give numerical results for the σ_2 problem, i.e. for $k = 2, n = 3$ using the subharmonicity preserving iterations. Although our theoretical results only cover smooth solutions, as indicated in the abstract and in the introduction, our methods appear able to handle non smooth solutions. The initial guess in all of our numerical experiments is taken as the finite difference approximation of the solution of the Poisson equation $\Delta u = 2\sqrt{f}$ in Ω with $u = g$ on $\partial\Omega$.

| | h | | | | |
|-------|------------------------|------------------------|------------------------|------------------------|------------------------|
| | $1/2^1$ | $1/2^2$ | $1/2^3$ | $1/2^4$ | $1/2^5$ |
| Error | $6.2328 \cdot 10^{-2}$ | $2.6556 \cdot 10^{-2}$ | $7.7836 \cdot 10^{-3}$ | $2.0616 \cdot 10^{-3}$ | $5.2449 \cdot 10^{-4}$ |
| Rate | | 1.23 | 1.77 | 1.92 | 1.97 |

TABLE 1. Maximum error with Test 1.

We use the following test functions on the unit cube $[0, 1]^3$:

Test 1: A smooth solution which is strictly convex, $u(x, y, z) = e^{x^2+y^2+z^2}$ so that $f(x, y, z) = 4(3 + x^2 + y^2 + z^2)e^{2(x^2+y^2+z^2)}$ and $g(x, y, z) = e^{x^2+y^2+z^2}$ on $\partial\Omega$.

Test 2: A smooth solution which is 2-convex but not convex. It is known that for a radial function $u(x) = \phi(r)$, $r = |x|$, $x \in \mathbb{R}^n$ the eigenvalues of D^2u are given by $\lambda_1 = \phi''(r)$ with multiplicity 1 and $\lambda_2 = \phi'(r)/r$ with multiplicity $n - 1$. See for example [7, Lemma 2.1]. It follows that with $u(x, y, z) = \ln(a + x^2 + y^2 + z^2)$, we have $\phi(r) = \ln(a + r^2)$ and we get $\Delta u = \frac{6a+2r^2}{(a+r^2)^2} \geq 0$, $S_2(D^2u) = 4\frac{3a-r^2}{(a+r^2)^3} \geq 0$, $\det D^2u = 2\frac{a-r^2}{(a+r^2)^2}$, in $[0, 1]^3$. With $a = 2$, $\det D^2u$ takes negative values in $[0, 1]^3$.

Test 3: A solution not in $H^2(\Omega)$, $u(x, y, z) = -\sqrt{3 - x^2 - y^2 - z^2}$ so that $f(x, y, z) = -(x^2 + y^2 + z^2 - 9)/(-3 + x^2 + y^2 + z^2)^2$ and $g(x, y, z) = -\sqrt{3 - x^2 - y^2 - z^2}$ on $\partial\Omega$.

Test 4: No exact solution is known. Here $f(x, y, z) = 1$ and $g(x, y, z) = 0$.

Test 5: A degenerate three dimensional Monge-Ampère equation. We take $f(x, y, z) = 0$ and $g(x, y, z) = |x - 1/2|$. We use the double iterative method based on (1.5).

Numerically, the solution computed may not satisfy $S^2D^2u^m \geq 0$. At those points we set both $S_2(D^2u^m)$ and $\det D^2u^m$ to 0 in (1.5). If the numerical value of $S_2(D^2u^m)$ is negative, then 0 is a better approximate value. Since $S_2(D^2u^m)$ is computed from u^m , the numerical value of $\det D^2u^m$ would also be inaccurate. Since u^m is expected to be an approximate solution of u for which $\det D^2u \geq 0$, a better approximation of $\det D^2u^m$ at any stage where the latter is negative is also 0. It would be interesting to analyze the effect of these rounding off errors on the overall numerical convergence of the method. For example, one may analyze the convergence of the inexact double iteration. Similar situations appear with inexact Newton's methods and inexact Uzawa algorithms.

The right hand side $f(x, y, z)$ can be computed from the exact solution $u(x, y, z)$ using the definition of $S_2(D^2u)$ as the sum of the 2×2 principal minors.

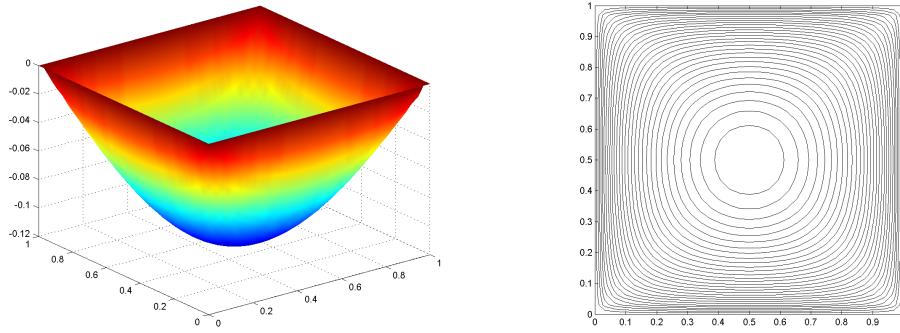
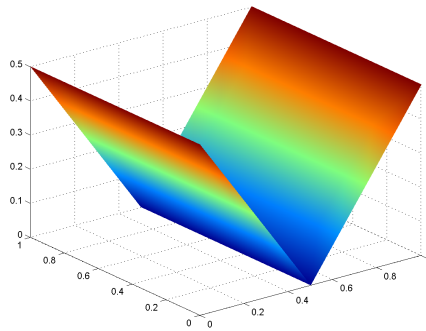
For all tests but Test 3, we used the direct solver (5.1). For Test 3, the Gauss-Seidel method was used since we needed values for small values of h . As expected, we have quadratic convergence (as $h \rightarrow 0$) for the smooth solutions of Tests 1 and 2 while enough data is not available to give the convergence rate for the singular solution of Test 3.

| | h | | | | |
|-------|------------------------|------------------------|------------------------|------------------------|------------------------|
| | $1/2^1$ | $1/2^2$ | $1/2^3$ | $1/2^4$ | $1/2^5$ |
| Error | $6.5241 \cdot 10^{-4}$ | $5.0653 \cdot 10^{-4}$ | $1.3850 \cdot 10^{-4}$ | $3.5587 \cdot 10^{-5}$ | $9.1276 \cdot 10^{-6}$ |
| Rate | | 0.36 | 1.87 | 1.96 | 1.96 |

TABLE 2. Maximum error with Test 2.

| | h | | |
|-------|------------------------|------------------------|------------------------|
| | $1/2^4$ | $1/2^5$ | $1/2^6$ |
| Error | $1.1084 \cdot 10^{-3}$ | $9.7971 \cdot 10^{-4}$ | $7.6618 \cdot 10^{-4}$ |
| Rate | | 0.18 | 0.35 |

TABLE 3. Maximum error with Test 3.

FIGURE 1. Test 4, $h = 1/2^5$. Graph and contour in plane $z = 1/2$.FIGURE 2. Test 5, $h = 1/2^4$. Graph in the plane $z = 1/2$.

In [3], it was argued based on numerical evidence that the Gauss-Seidel method (5.4) is faster than a certain variant of the direct solver (5.1) for singular solutions. In our implementation we saw evidence of the contrary, that is, the Gauss-Seidel method is

less efficient. We note that the Gauss-Seidel method requires much more loops which are not efficient in MATLAB.

7. CONCLUDING REMARKS

Remark 7.1. *Although the pseudo-transient and time marching methods introduced in [2] work as well for k -Hessian equations, and apply to more general fully nonlinear equations, the subharmonicity preserving iterative methods introduced in this paper are parameter free. All these type of methods can be accelerated with fast Poisson solvers and multigrid methods.*

Remark 7.2. *When it comes to numerical methods for fully nonlinear equations, there are two types of convergence to study. Since the equations are nonlinear, they must be solved iteratively. One must then address the convergence to the discrete solution of the iterative methods used. The second type of convergence is the convergence of the numerical solution to the exact solution as the discretization parameter converges to 0. We have addressed both types of convergence in this paper.*

Remark 7.3. *Existence of a discrete solution and convergence (as the mesh size $h \rightarrow 0$), for finite difference discretization of smooth solutions of fully nonlinear equations, are not often discussed. It is clear that convergence does not simply follow from the consistency of standard finite difference discretization of the second order derivatives. For viscosity solutions, convergence of monotone, stable and consistent schemes follows immediately from the theory of Barles and Souganidis.*

Remark 7.4. *The iterative method (1.3) can be viewed as a linearization of the fully nonlinear equation (1.1). It is possible to linearize (1.1) in ways different from (1.2) and (1.3). See for example the methods described in [2]. The iterative method (1.3) has been shown numerically to select discrete solutions which converge to non smooth solutions. Since (1.3) consists of a sequence of Poisson equations, the numerical solution of (1.1) can now be tackled with any good numerical method.*

REFERENCES

- [1] Awanou, G.: On standard finite difference discretizations of the elliptic Monge-Ampère equation (2014). <http://homepages.math.uic.edu/~awanou/up.html>
- [2] Awanou, G.: Pseudo transient continuation and time marching methods for Monge-Ampère type equations (2014). <http://arxiv.org/abs/1301.5891>. To appear in Advances in Comp. Math.
- [3] Benamou, J.D., Froese, B.D., Oberman, A.M.: Two numerical methods for the elliptic Monge-Ampère equation. *M2AN Math. Model. Numer. Anal.* **44**(4), 737–758 (2010)
- [4] Chen, Y.: Efficient and robust solvers for Monge-Ampère equations. Ph.D. thesis, Clarkson University (2010)
- [5] Chen, Y., Fulton, S.R.: An adaptive continuation-multigrid method for the balanced vortex model. *J. Comput. Phys.* **229**(6), 2236–2248 (2010)
- [6] Chou, K.S., Wang, X.J.: A variational theory of the Hessian equation. *Comm. Pure Appl. Math.* **54**(9), 1029–1064 (2001)
- [7] Felmer, P.L., Quaas, A.: On critical exponents for the Pucci’s extremal operators. *Ann. Inst. H. Poincaré Anal. Non Linéaire* **20**(5), 843–865 (2003)
- [8] Feng, X., Neilan, M.: Analysis of Galerkin methods for the fully nonlinear Monge-Ampère equation. *J. Sci. Comput.* **47**(3), 303–327 (2011)

- [9] Gavitone, N.: Hessian equations, quermassintegrals and symmetrization. Ph.D. thesis, Università di Napoli Federico II (2009)
- [10] Gilbarg, D., Trudinger, N.S.: Elliptic partial differential equations of second order. Classics in Mathematics. Springer-Verlag, Berlin (2001). Reprint of the 1998 edition
- [11] Headrick, M., Wiseman, T.: Numerical Ricci-flat metrics on K3. *Classical and Quantum Gravity* **22**(23), 4931–4960 (2005)
- [12] Hoffman, A.J., Wielandt, H.W.: The variation of the spectrum of a normal matrix. *Duke Math. J.* **20**, 37–39 (1953)
- [13] Horn, R.A., Johnson, C.R.: Matrix analysis. Cambridge University Press, Cambridge (1990). Corrected reprint of the 1985 original
- [14] Johnson, C.G.L.: Estimates near plane portions of the boundary for discrete elliptic boundary problems. *Math. Comp.* **28**, 909–935 (1974)
- [15] Kuo, H.J., Trudinger, N.S.: Linear elliptic difference inequalities with random coefficients. *Math. Comp.* **55**(191), 37–53 (1990)
- [16] Lieberman, G.M.: Second order parabolic differential equations. World Scientific Publishing Co. Inc., River Edge, NJ (1996)
- [17] Oberman, A.M.: Numerical methods for the 2-Hessian partial differential equation. Submitted, 2015
- [18] Sorensen, D.C., Glowinski, R.: A quadratically constrained minimization problem arising from PDE of Monge-Ampère type. *Numer. Algorithms* **53**(1), 53–66 (2010)
- [19] Thomée, V.: Discrete interior Schauder estimates for elliptic difference operators. *SIAM J. Numer. Anal.* **5**, 626–645 (1968)
- [20] Wang, X.J.: The k -Hessian equation. In: Geometric analysis and PDEs, *Lecture Notes in Math.*, vol. 1977, pp. 177–252. Springer, Dordrecht (2009)

DEPARTMENT OF MATHEMATICS, STATISTICS, AND COMPUTER SCIENCE, M/C 249. UNIVERSITY OF ILLINOIS AT CHICAGO, CHICAGO, IL 60607-7045, USA

E-mail address: `awanou@uic.edu`

URL: `http://www.math.uic.edu/~awanou`